# Introduction to Rice HPCToolkit on Early Access BlueGene/Q

Mark W. Krentel

Department of Computer Science

Rice University

krentel@rice.edu

**http://hpctoolkit.org**

Saturday, April 28, 2012

# HPCToolkit Basic Features

- **Run application natively, every 100-1,000 times per second, interrupt program and record snapshot of call stack.**

- **Combine sampling data with binary analysis of program structure: loops, inline functions, etc.**

- **Present top-down, bottom-up and flat views of calling context tree (CCT) and time-sequence trace view.  Costs are displayed per source line in the context of their call path.**

- **Can sample on Wallclock (itimer) and Hardware Performance Counter Events (PAPI preset and native events).**

# Advantages of Sampling

- **Run application natively at full optimization.**

- **Analyze program binary, no changes to source code.**

- **Low overhead, typically < 5%, overhead is proportional to sampling rate, not number of function calls.**

Saturday, April 28, 2012

# HPCToolkit Advanced Features

- **Finely-tuned unwinder to handle multi-lingual, fully-optimized code, no frame pointers, broken return pointers, stack trolling, etc.**

- **Derived metrics -- compute flops per cycle, or flops per memory reads, etc. and attribute to lines in source code.**

- **Compute strong and weak scaling loss, for example:**

  **strong:  8 * (time at 8K cores) - (time at 1K cores)**

  **weak:    (time at 8K cores and 8x size) - (time at 1K cores)**

- **Blame shifting -- when thread is idle or waiting on a lock, blame the working threads or holder of lock.**

- **Load imbalance -- display distribution and variance in metrics across cores and threads.**

4

# Getting Started with HPCToolkit

- **Add to PATH:**

  **/home/projects/hpc/pkgs/hpctoolkit/bin**

- **Compile source files natively with full optimization, add -g to CFLAGS  (for source lines).**

- **Use hpclink to link application with hpctoolkit code.**

  **hpclink  mpicc  -o myprog  file.o ...  -llib ...**

- **Launch program with HPCRUN environ variables.**

  **HPCRUN_EVENT_LIST='PAPI_TOT_CYC@15000000,**
  **PAPI_FP_OPS@1000000'**
  **HPCRUN_TRACE=1    (for tracing)**
  **qsub -A <project> -t <time> -n <nodes>  ...  \**
  **--env HPCRUN_EVENT_LIST='...':HPCRUN_TRACE=1  \**
  **myprog arg  ...**

5

# Getting Started, cont'd.

- **Use hpcstruct to analyze program binary.**

  hpcstruct  myprog

  =>  myprog.hpcstruct

- **Use hpcprof or hpcprof-mpi to combine .hpcstruct file with measurements directory (use '+' for subdirectories).**

  hpcprof  -S myprog.hpcstruct  \

    -I /path/to/myprog/source/tree/+  \

    hpctoolkit-myprog-measurements-jobid

  ==>  hpctoolkit-myprog-database-jobid

- **Use hpcviewer and hpctraceview (if enabled tracing) to view results.**

  hpcviewer  hpctoolkit-myprog-database-jobid

  hpctraceviewer  hpctoolkit-myprog-database-jobid

6

# Where to find HPCToolkit

- **Home page:**

  **http://hpctoolkit.org/**

- **On veas:**

  **/home/projects/hpc/pkgs/hpctoolkit/bin**

- **Source code available for anonymous svn checkout at the SciDAC Outreach Center (hpctoolkit project).**

  **https://outreach.scidac.gov/projects/hpctoolkit/**

- **Prebuilt versions of the viewer and traceviewer also available at the SciDAC Outreach Center (hpcviewer project).**

  **https://outreach.scidac.gov/projects/hpcviewer/**

# HPCToolkit Capabilities at a Glance



Attribute Costs to Code

Pinpoint & Quantify Scaling Bottlenecks

Assess Imbalance and Variability

Analyze Behavior over Time

Shift Blame from Symptoms to Causes

Associate Costs with Data

RICE

**hpctoolkit.org**

# Call Path Profiling

**Measure and attribute costs in context**

**sample timer or hardware counter overflows**

**gather calling context using stack unwinding**



Call path sample

Calling context tree

- return address
- return address
- return address
- instruction pointer

**Overhead proportional to sampling frequency...**
**...not call frequency**

# Understanding Temporal Behavior

- **Profiling compresses out the temporal dimension**
  - temporal patterns, e.g. serialization, are invisible in profiles

- **What can we do? Trace call path samples**
  - sketch:
    - N times per second, take a call path sample of each thread
    - organize the samples for each thread along a time line
    - view how the execution evolves left to right
    - what do we view?
      assign each procedure a color; view a depth slice of an execution

Processes

Time

Call stack

# AMG2006: 8PE x 8 OMP Threads



OpenMP loop in `hypre_BoomerAMGRelax` using static scheduling has load imbalance; threads idle for a significant fraction of their time

# Code-centric view: `hypre_BoomerAMGRelax`

# Serial Code in AMG2006 8 PE, 8 Threads



7 worker threads are idle in each process while its main MPI thread is working

# Pinpointing and Quantifying Scalability Bottlenecks
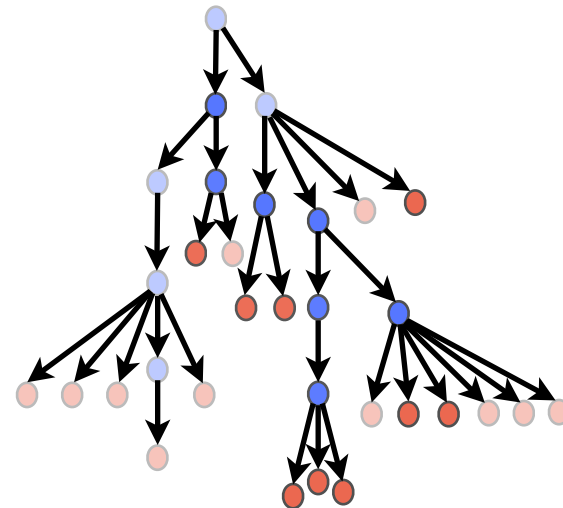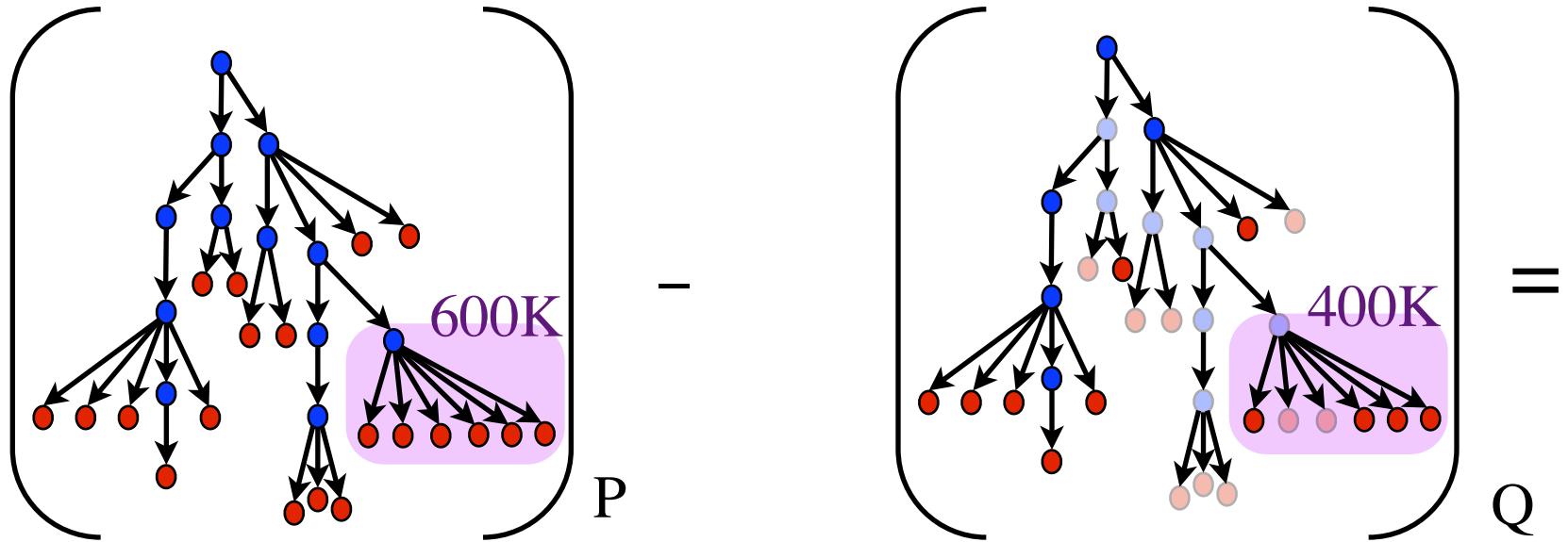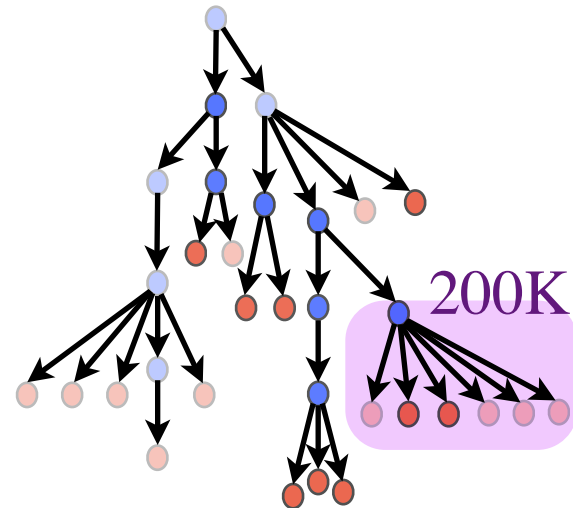


$$P \quad - \quad Q \quad =$$

400K

# Pinpointing and Quantifying Scalability Bottlenecks
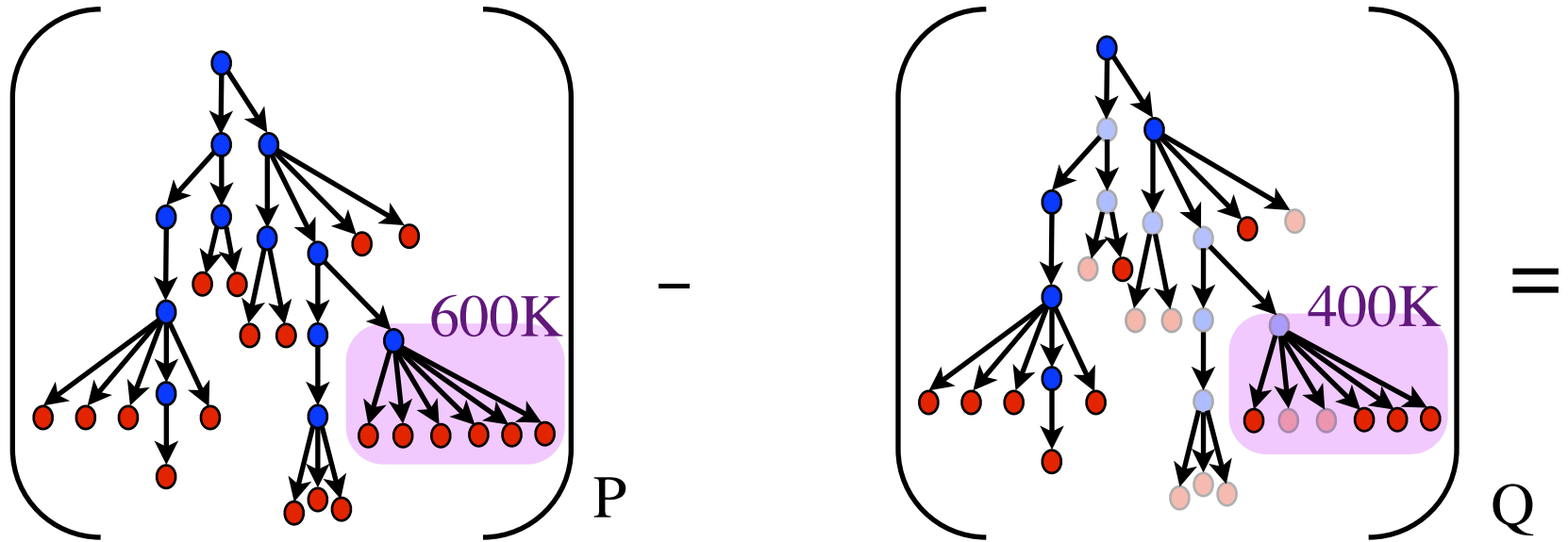
# Pinpointing and Quantifying Scalability Bottlenecks



$P \times$ [ tree diagram 600K ] $P$ $-$ $Q \times$ [ tree diagram 400K ] $Q$ $=$ [ tree diagram 200K ]

coefficients for analysis of strong scaling